

The NIH Comparative Genomics Resource (CGR)

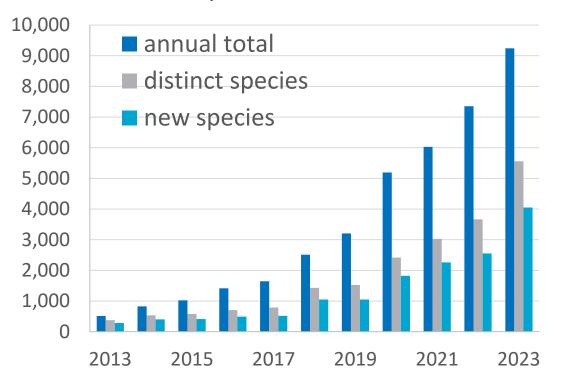
Terence D. Murphy – CGR Project Lead



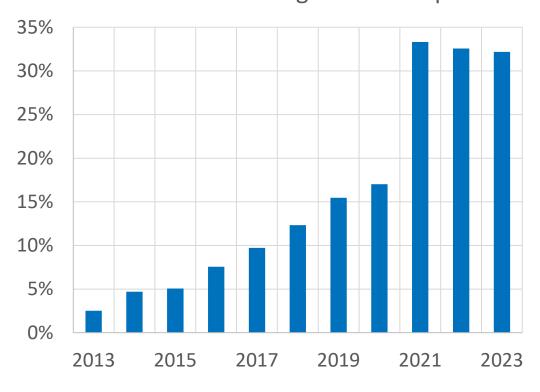
Look at all the genomes!

Current eukaryotic genomes (as of Dec 31 2023): 36,593 genomes 15,453 species





Percent with ContigN50 > 1 Mbp







Why CGR?

Comparative genomics investigators face several limitations and challenges

- Multiple different user interfaces
- Limited number of organisms supported
- Siloed data and applications
- Must download data to apply tools
- Limited scalability
- Data quality issues

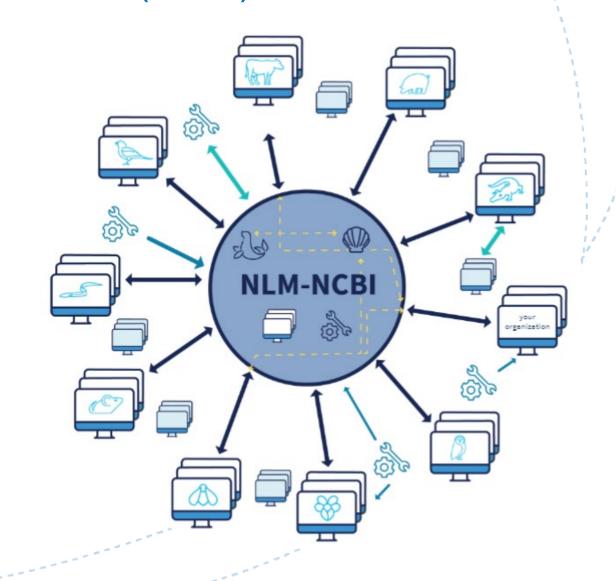
NCBI

Solution

NIH Comparative Genomics Resource (CGR)

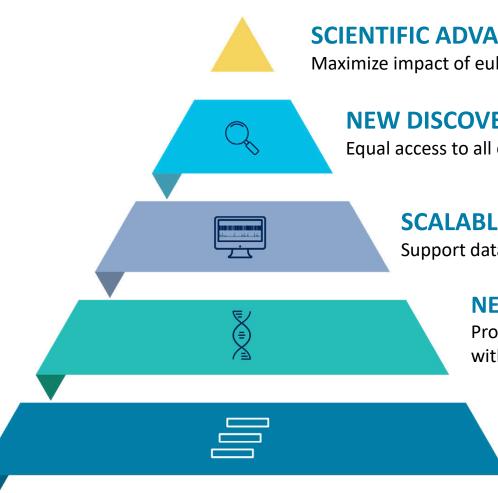
What: CGR maximizes the impact of eukaryotic research organisms and their genomic data to biomedical research.

How: CGR facilitates reliable comparative genomics analyses through community collaboration and an NCBI genomics toolkit. The toolkit includes high-quality data, tools, and interfaces for connecting community-provided resources with NCBI.





CGR Benefits



SCIENTIFIC ADVANCEMENT

Maximize impact of eukaryotic research organisms to biomedical research

NEW DISCOVERY AMPLIFICATION

Equal access to all eukaryotic organism data with better connections to community resources

SCALABLE ANALYSES

Support data growth with emerging big data approaches

NEW AND IMPROVED COMPARATIVE GENOMICS TOOLS

Promote high-quality data submission, exploration, analysis, and retrieval with seamless user experiences

HIGH QUALITY GENOMIC DATA

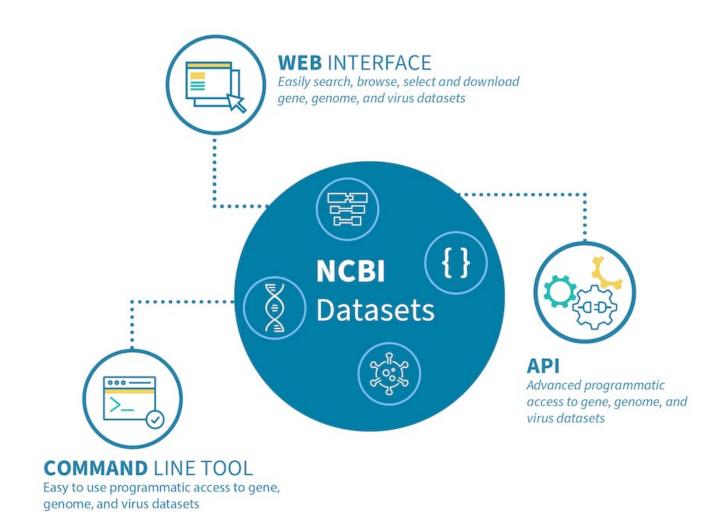
Serve standardized, uncontaminated, and consistently annotated eukaryotic genomic data from NCBI Archives





NCBI Datasets

- Easy-to-use web and programmatic interfaces
- Data delivery in formats that support both web and programmatic users
- Helpful documentation and tutorials









NCBI Datasets Genome Data Package

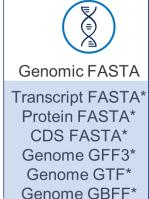
Data

- Genome FASTA
- Annotation (GFF, GTF, GBFF)
- Transcript, Protein,
 CDS FASTA

GENOME



GENOME DATA PACKAGE





Data report Sequence report Dataset catalog

Metadata

- Contains metadata
 from Assembly, Gene,
 BioSample, and
 BioProject databases
- JSON lines format

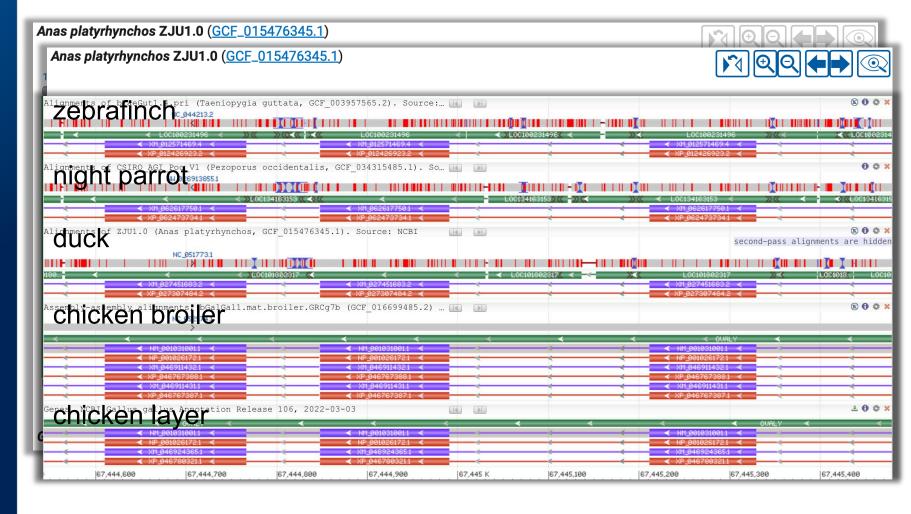




Comparative Genome Viewer (CGV)

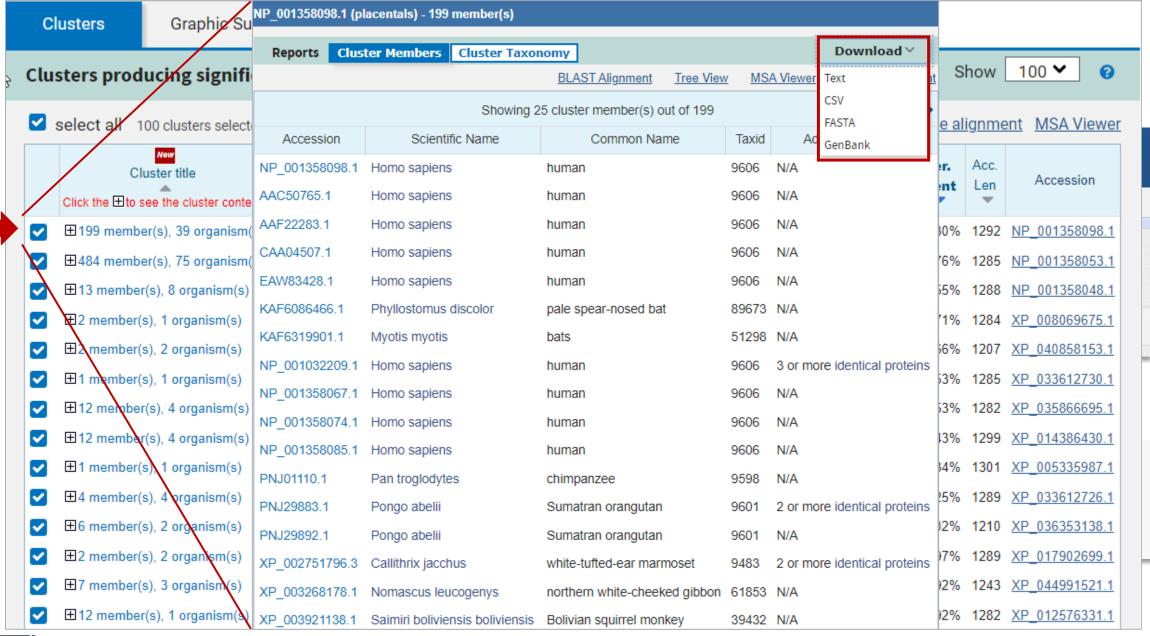
- Compare genomes through assembly alignments
- Over 700 alignments for over 300 species!
- Preprint: PMC10705539

https://www.ncbi.nlm.nih.gov/genome/cgv/



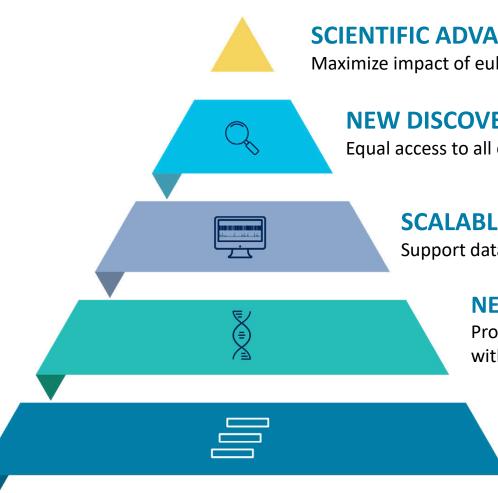


https://blast.ncbi.nlm.nih.gov/Blast.cgi



National Library of Medicine National Center for Biotechnology Information

CGR Benefits



SCIENTIFIC ADVANCEMENT

Maximize impact of eukaryotic research organisms to biomedical research

NEW DISCOVERY AMPLIFICATION

Equal access to all eukaryotic organism data with better connections to community resources

SCALABLE ANALYSES

Support data growth with emerging big data approaches

NEW AND IMPROVED COMPARATIVE GENOMICS TOOLS

Promote high-quality data submission, exploration, analysis, and retrieval with seamless user experiences

HIGH QUALITY GENOMIC DATA

Serve standardized, uncontaminated, and consistently annotated eukaryotic genomic data from NCBI Archives



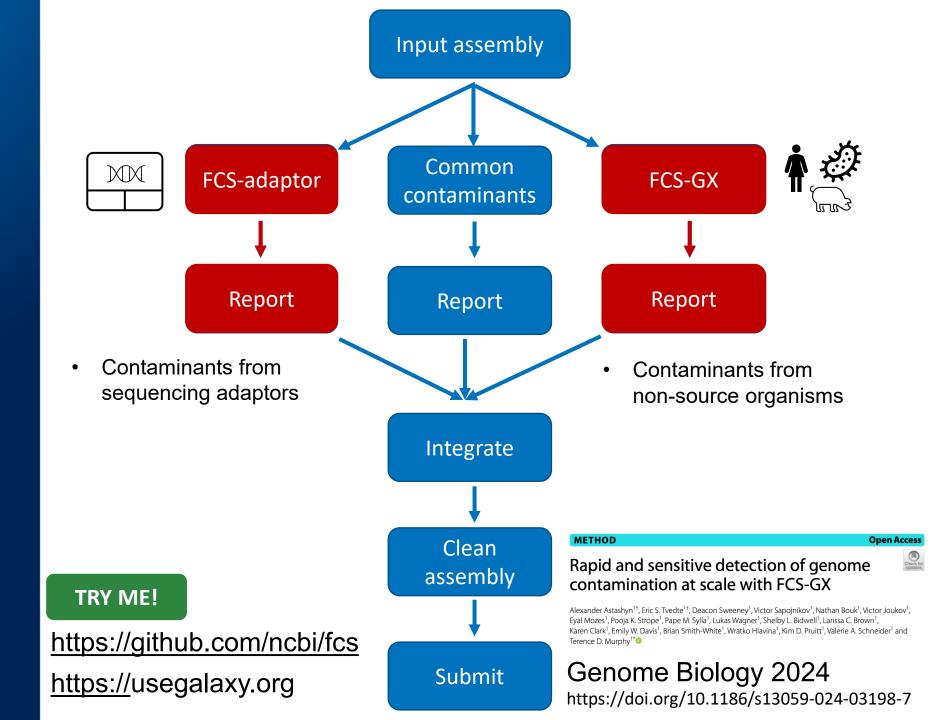


Foreign Contaminant Screening (FCS)

CGR Goals:

- Modular
- Easy
- Reliable
- High sensitivity and specificity
- Fast
- Same code used for GenBank submission screening





FCS-GX identifies contamination across the tree of life

36.8 Gbp of suspected contamination

- 0.16% bases
- 1.30% sequences
- Nearly 3 million annotated proteins
- Half from 161 assemblies

GenBank genome type

			<u> </u>						
contaminant type			4	*	***	1%	TOTAL		
	Metazoa	1.11	0.09	1.27	0.92	0.18	3.56	perce	ent contamina by length
	Fungi	1.44	0.35	0.70	0.56	0.10	3.14		0.001%
	Viridiplantae 🏅	0.42	0.07	0.01	0.05	0.18	0.74		0.01%
	Other Eukaryotes	2.09	0.01	0.06	0.03	0.43	2.62		0.1%
	Prokaryotes 💦	15.41	1.97	2.58	1.54	4.98	26.49		
	Viruses 🕌	0.14	0.01	0.02	0.01	0.02	0.20		
	Synthetic ‡	0.01	<0.01	0.02	<0.01	0.01	0.05		
	TOTAL	20.63	2.51	4.66	3.10	5.90	36.80		

values in Gbp

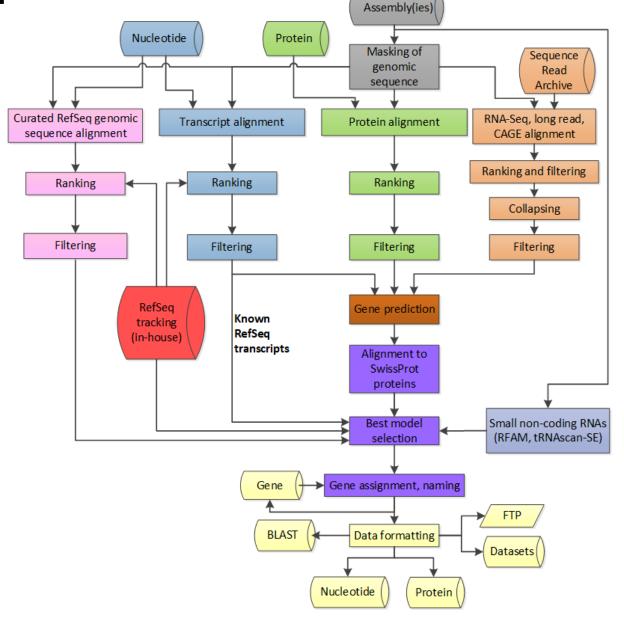
EGAP: Eukaryotic Genoma Annotation Pipeline

Used by NCBI to annotate ~1100 species Evidence used for gene prediction:

- ✓ ESTs
- ✓ cDNAs
- ✓ Same and cross-species proteins
- ✓ RNA-Seq
- ✓ PacBio IsoSeq, ONT transcriptomes
- ✓ CAGE
- ✓ PhyloCSF (for internal review)

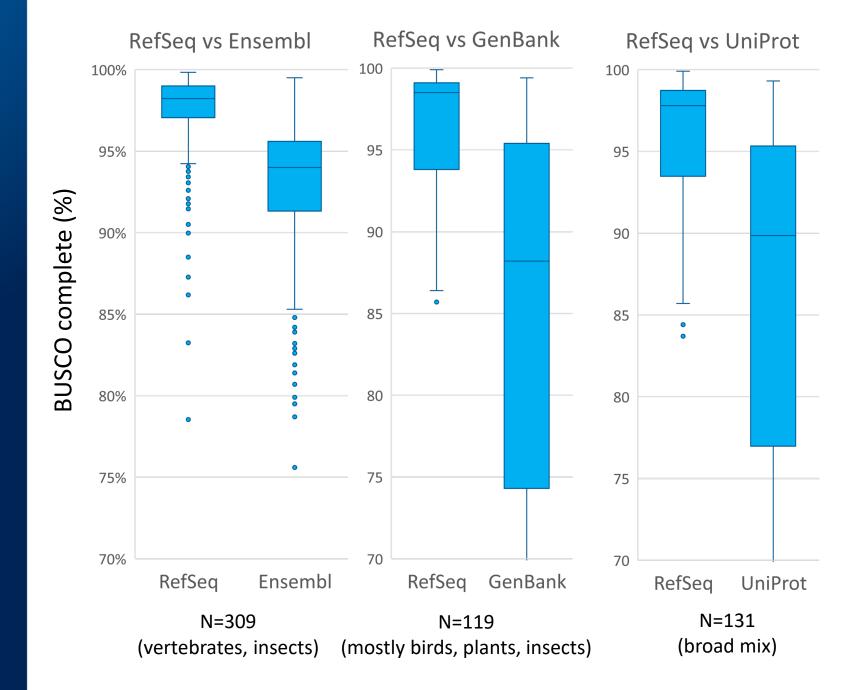
Outputs:

- ✓ Protein-coding genes
- ✓ Non-coding genes and pseudogenes
- ✓ Orthologs
- ✓ Gene Ontology
- ✓ Expression



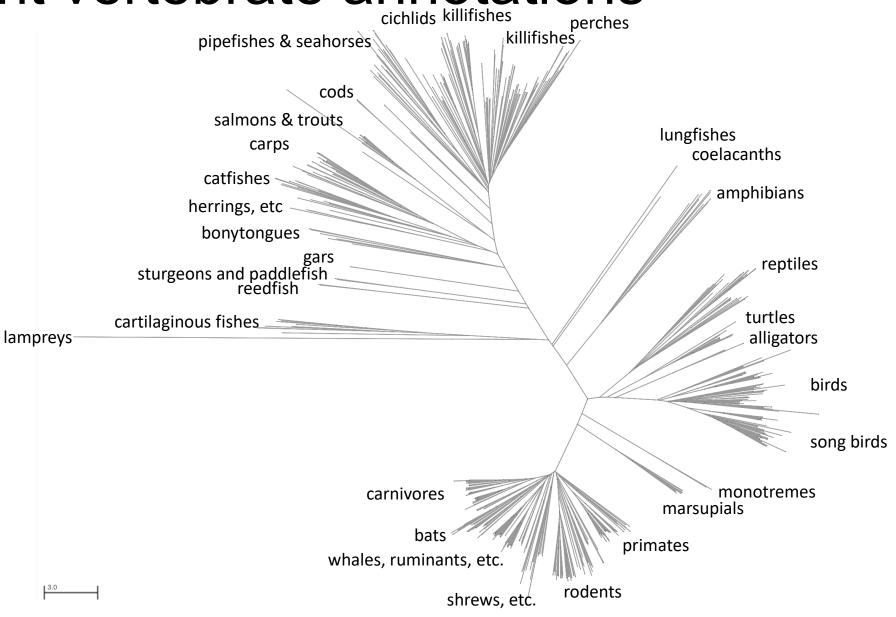


More complete
BUSCO
models found
in RefSeq
annotations





>580 current vertebrate annotations

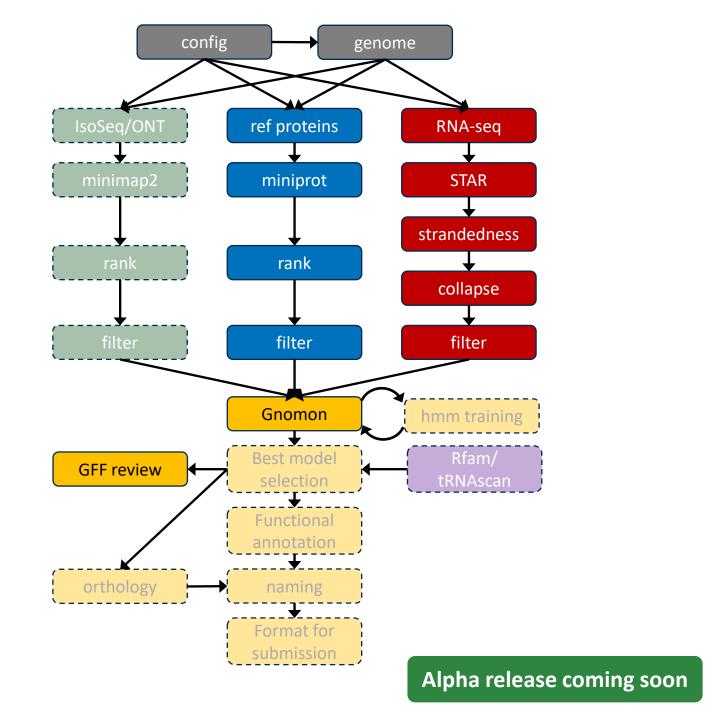




Public EGAPx

Goals:

- NCBI EGAP that you can run yourself
- Nextflow pipeline available using Docker container
- Produce submissionready annotation
- Work with public or private data
- Cloud-ready



Expanding annotation within a species

- NCBI RefSeq
 - Key organisms based on users or taxonomic placement
 - Focus on "reference" genomes (typically 1 per species)
 - Aim for highest quality annotation
- Additional assemblies
 - Annotate through EGAPx or LiftOff-type approach
 - Promote submission with annotation into GenBank archive
- Future support for pan-genomes combining the above



How Do I Get Involved?



Reach out to us at cgr@nlm.nih.gov



Visit the CGR website ncbi.nlm.nih.gov/cgr and click the yellow Feedback button on the bottom right of the page

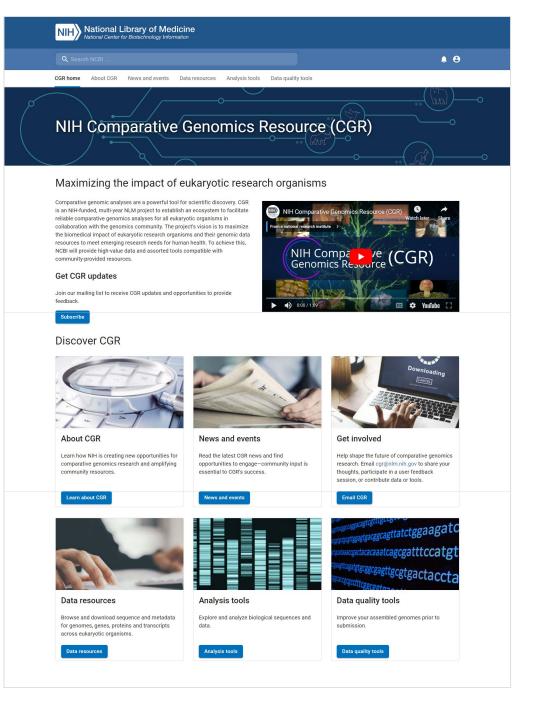


Sign up for our mailing list - bit.ly/Subscribe CGR



Look out for future meetings, workshops, webinars, surveys, small group sessions, user testing, and interviews to inform the development process





Thank you

This research was supported by the Intramural Research Program of the NIH, National Library of Medicine.

CGR product team leads & admin

Terence Murphy

Nuala O'Leary

Sanjida Rangwala

Tom Madden

Aron Marchler-Bauer

Anne Ketter

Katya Sukharnikov

NCBI Leadership

Francoise Thibaud-Nissen

Valerie Schnieder

Kim Pruitt

Steve Sherry

Developers, subject matter experts and support

Kelsey Aadland	Emily Davis	Avi Kimchi	Dong-Ha Oh	Mirian Tsuchiya
Victor Ananiev	Olga Ermolaeva	Vamsi Kodali	Marina Omelchenko	Eric Tvedte
Preye Akuiyibo	Vladislav Evgeniev	Raymond Koehler	Yan Raytselis	Joel Virothaisakun
Ray Anderson	Robert Falk	Christopher Lanczycki	Dmitry Rudnev	Lukas Wagner
Alex Astashyn	Amelia Fong	Vadim Lotov	Victor Sapojnikov	Craig Wallin
Andrea Asztalos	Steven Gaudaen	Shennan Lu	Robert Smith	Jiyao Wang
Hena Bajwa	Marc Gwadz	Tom Madej	Alexandre Souvorov	Mingzhang Yang
Greg Boratyn	Vichet Hem	Gabi Marchler	Deacon Sweeney	Jian Ye
Evgeny Borodin	Wratko Hlavina	Patrick Masterson	Pooja Strope	Irena Zaretskyna
Nathan Bouk	Jinna Hoffman	Scott McGinnis	Pape Sylla	Dachuane Zhang
Christiam Camacho	Brad Holmes	Yuri Merezhuk	Jovany Tinne	
Eric Cox	Victor Joukov	Eyal Mozes	John Torcivia	

Contact us info@ncbi.nlm.nih.gov



Watch NCBI News for updates!

http://www.ncbi.nlm.nih.gov/news/
https://www.youtube.com/user/NCBINLM